

Representation:

$x^{(t)}, y^{(t)}$ → t^{th} element in sequence
 T_x, T_y = length of sequence
 $x^{(i)(t)}$ → for i^{th} training example
 $T_x^{(i)}, T_y^{(i)}$ → length of i^{th} training sequence
can change per training example

Vocabulary / dictionary: store all words

then use one hot representation having a vector per word full of zeros except i^{th} location where i is the position of word in dictionary

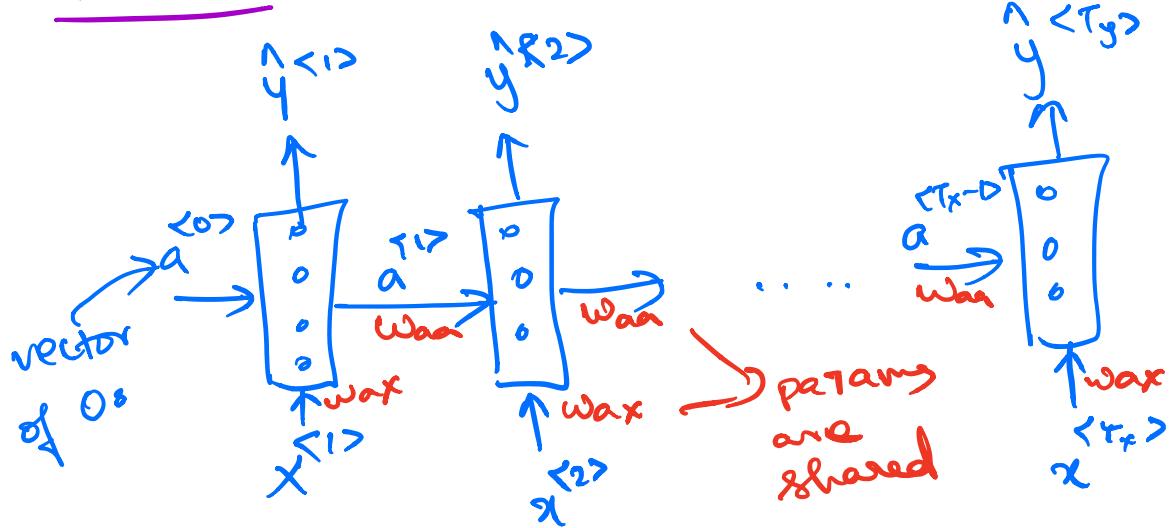
$x^{(t)}$ is a one hot vector.

$$x^{(t)} = \begin{cases} 1, & \text{for } i = \text{dict[i]}[\text{word}] \\ 0, & \text{elsewhere} \end{cases}$$

we cannot use standard NN:

- i) input and output can be of different length
- ii) doesn't share features learned from different parts of the text.

RNN:



But this RNN doesn't consider words after it. (ie) it processes sentences only from left to right

forward prop:

$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g(W_{aa}a^{<0>} + W_{ax}x^{<1>} + b_a)$$

tanh/relu
y^{<1>} = g(W_{ya}a^{<1>} + b_y)

$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$

$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

(100, 100) 100 (100, 10, 000)

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$

↑ ↑ ↑

$$a^{<t>} = g(W_a [a^{<t-1>} \mid x^{<t>}] + b_a)$$

$$\begin{bmatrix} 100 \\ W_{aa} \mid W_{ax} \\ 100 \end{bmatrix} = W_a$$

(100, 10100)

$$[a^{<t-1>} \mid x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$$

↑ 100 ↓ 10000
↓ 10000 ↑ 10100

$$[W_{aa} \mid W_{ax}] \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = W_{aa}a^{<t-1>} + W_{ax}x^{<t>}$$

Back prop:

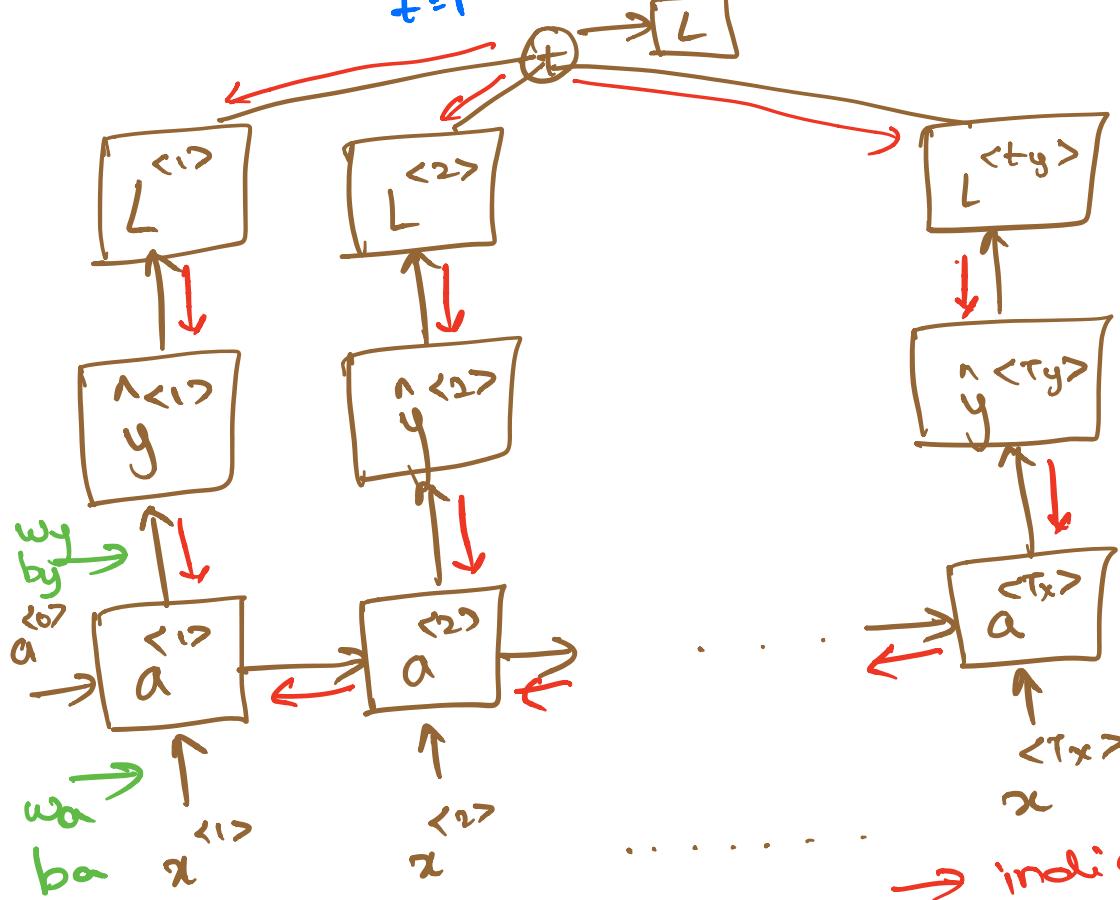
Back prop through time

cross entropy
loss
↓

$$L(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>}$$

$$- (1 - y^{<t>}) \log (1 - \hat{y}^{<t>})$$

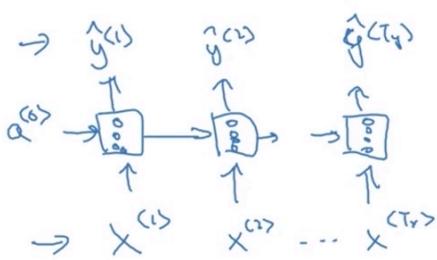
$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$



we can also have data where

$$T_x \neq T_y$$

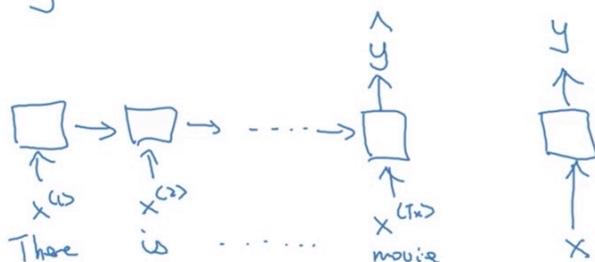
$$T_x = T_y$$



Many-to-many

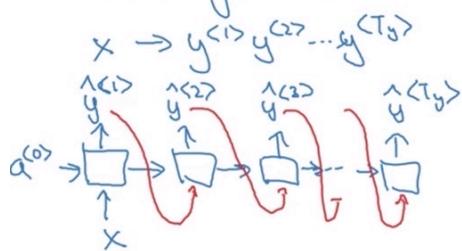
Sentiment classification

$$\begin{aligned} x &= \text{text} \\ y &= 0/1 \quad 1 \dots 5 \end{aligned}$$



One-to-one

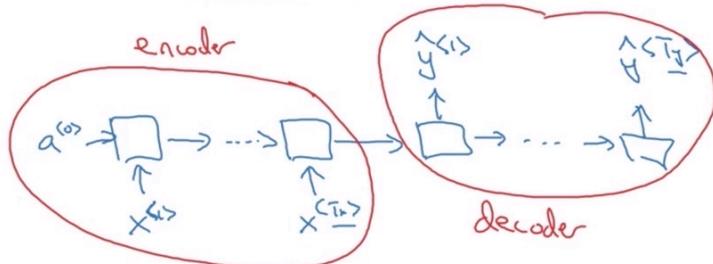
Music generation



One-to-many

$$x = \emptyset$$

Machine translation



Many-to-many ($T_x \neq T_y$)

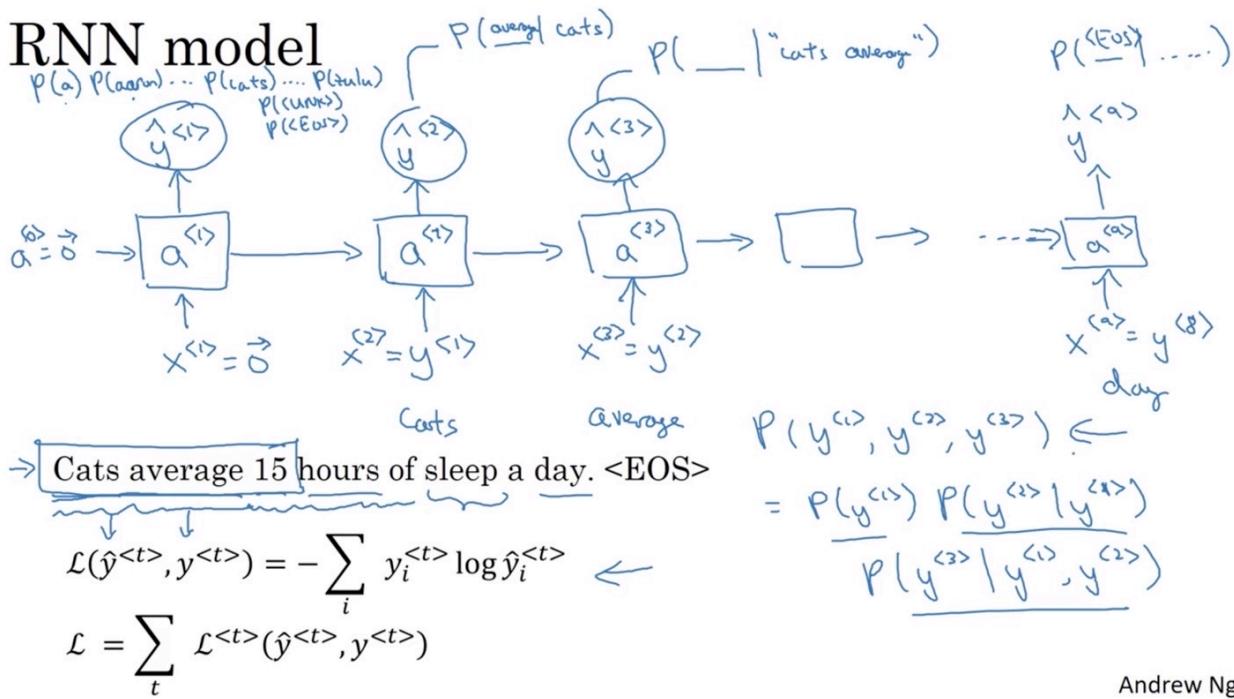
Language model:

sentence $\rightarrow p(\text{sentence})$

Training set: large corpus of english text.

input → tokenize (add $\langle \text{EOS} \rangle$)
 get $y^{<t>} (x^{<t>} = y^{<t-1>})$
 for unknown word we $\langle \text{UNK} \rangle$

RNN model



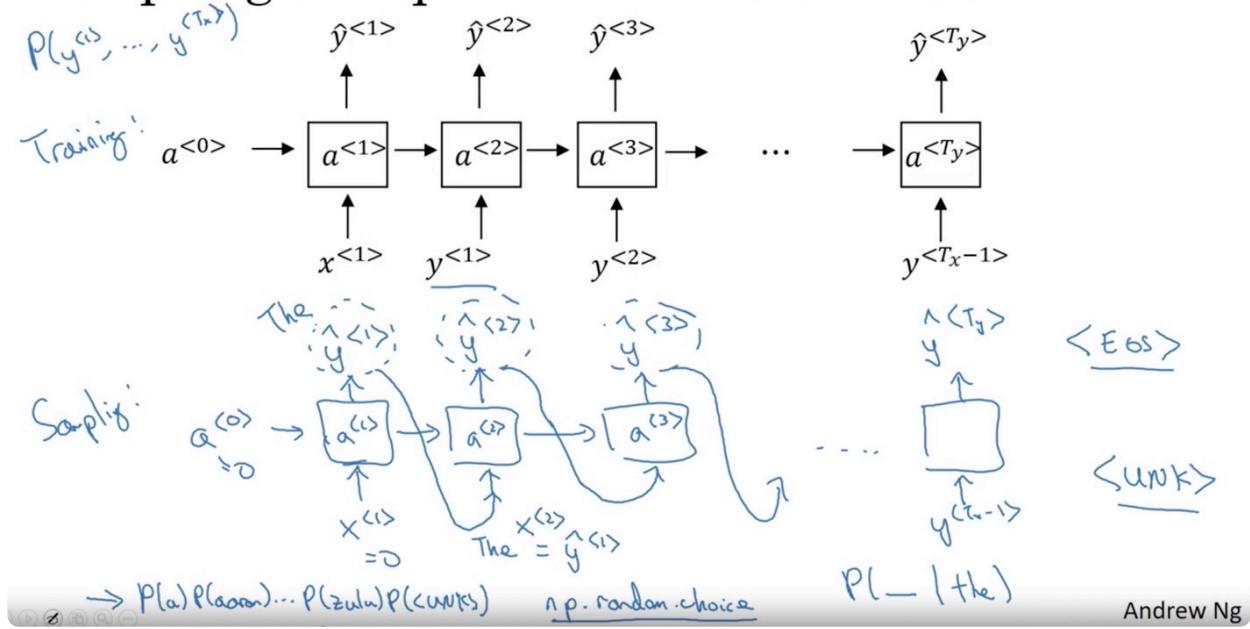
Andrew Ng

we can also sample sentences from this network using random sampling

To stop:

- i) sample 10-... words
- ii) add <EOS> to vocabulary
and stop when <EOS>
is sampled.

Sampling a sequence from a trained RNN



we can also build character level language model (instead of word level) in which case vocab will have alphabets instead of words

In character level \rightarrow we don't have to worry about unknown words

but we need longer sequences and hence performs worse when capturing long term dependencies.

In RNNs, we have the problem of vanishing gradients.

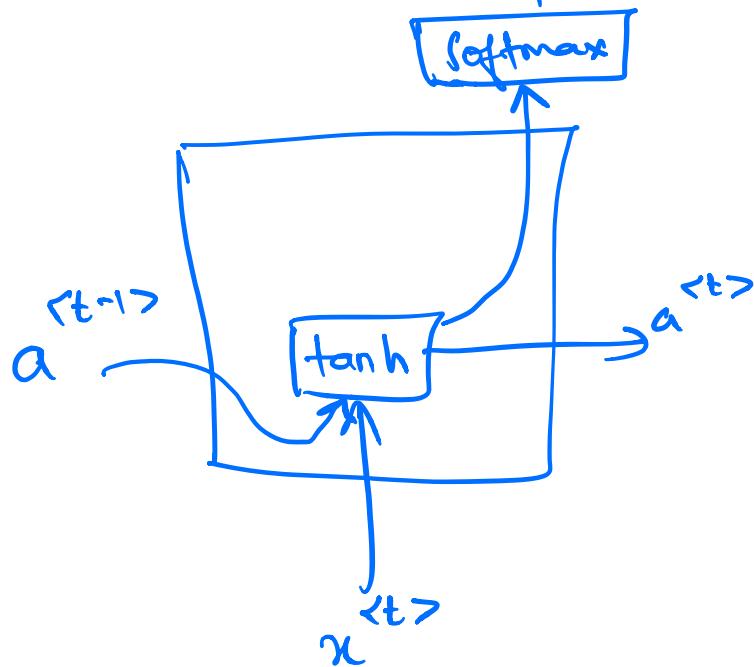
so, $y^{(100)}$ tends to be strongly influenced by few previous $y^{(t)}$ but not by $y^{(t)}$'s earlier in the sequence.

(\leftarrow) RNNs fail to capture long range dependencies.

for exploding gradients: apply gradient clipping.

Gated recurrent unit:

RNN unit:



GRU:

c = memory cell

$$c^{<t>} = a^{<t>}$$

candidate $\rightarrow c^{<t>} = \tanh(w_c [c^{<t-1>}, x^{<t>}] + b_c)$

to replace $c^{<t>}$

we have a gate: Γ_u (gamma)
 \approx (update gate)

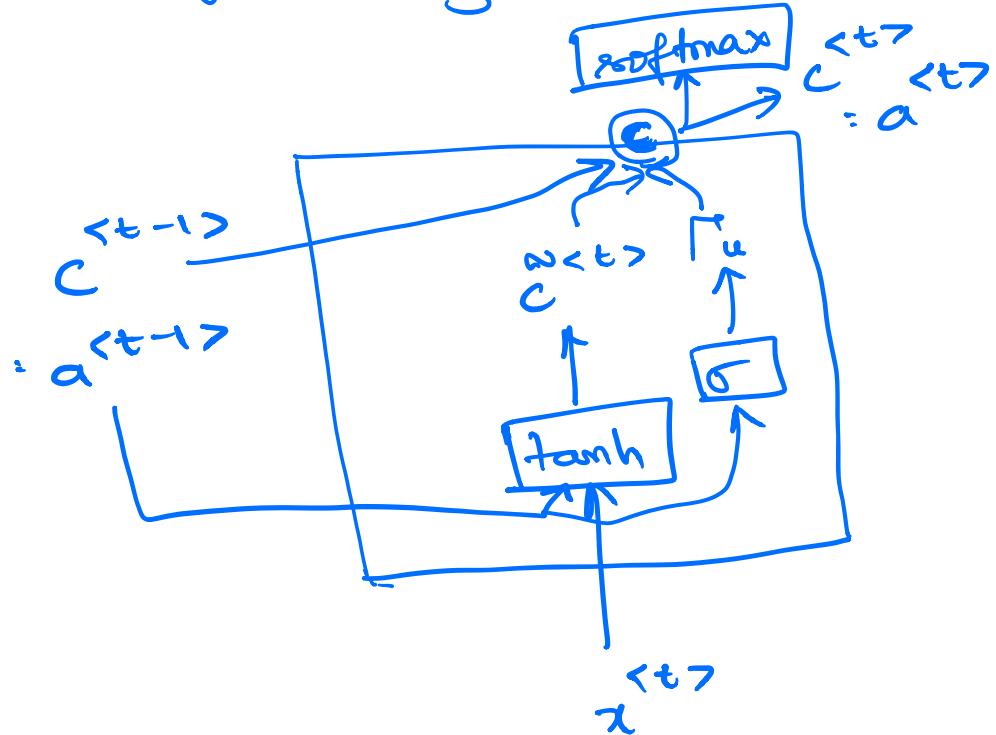
mostly 0/1 $\rightarrow \Gamma_u = \sigma(w_u [c^{<t-1>}, x^{<t>}] + b_u)$

when to update c + what to update

$$C^{<t>} = \Gamma_0 * C^{<t>} + (1 - \Gamma_0) * C^{<t-1>}$$

set $C^{<t>} = C^{<t>}$ when gate = 1

when gate = 0, $C^{<t>} = C^{<t-1>}$



as Γ_u will be close to zero,
 $C^{<t>} = C^{<t-1>}$, this GRU wont
 suffer with vanishing
 gradient

Full GRU:

$$\tilde{C}^{(t)} = \tanh(W_c [r_{\text{c}} * C^{(t-1)}, x^{(t)}] + b_c)$$

↑
relevance of
 $C^{(t-1)}$ to $C^{(t)}$

$$r_u = \sigma(W_u [C^{(t-1)}, x^{(t)}] + b_u)$$

$$r_r = \sigma(W_r [C^{(t-1)}, x^{(t)}] + b_r)$$

$$C^{(t)} = r_u * \tilde{C}^{(t)} + (1 - r_u) * C^{(t-1)}$$

$a^{(t)} = C^{(t)}$

Long short term memory: (LSTM)

no longer the case that
 $a^{(t)} = C^{(t)}$

$$\tilde{C}^{(t)} = \tanh(W_c [a^{(t-1)}, x^{(t)}] + b_c)$$

we have two update gate

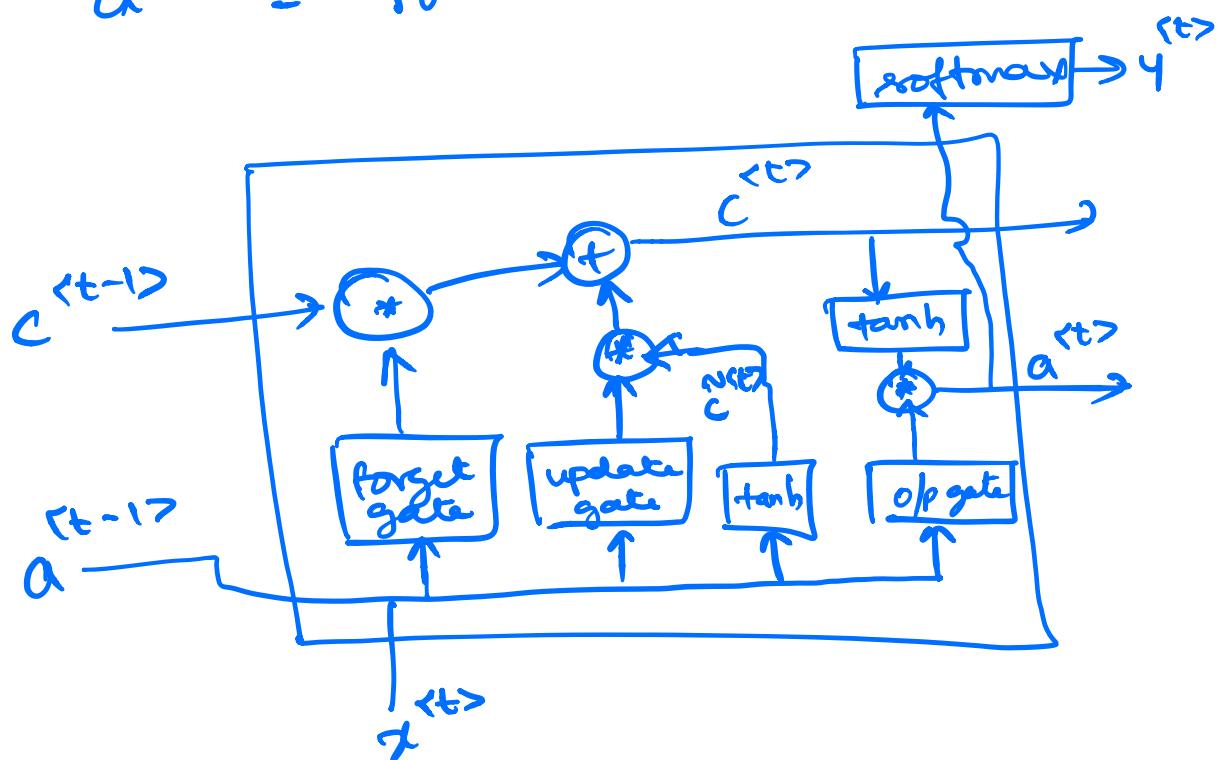
$$\text{update gate}_u = \sigma(w_u [a^{t-1}, x^{t}] + b_u)$$

$$f_f = \sigma(w_f [a^{t-1}, x^{t}] + b_f)$$

$$\text{output } f_o = \sigma(w_o [a^{t-1}, x^{t}] + b_o)$$

$$c^{t} = f_u * c^{t-1} + f_f * c^{t-1}$$

$$a^{t} = f_o * \tanh(c^t)$$



If update / forget gate works correctly, $c^{(100)}$ can be $c^{(1)}$.

so this helps in overcoming vanishing gradients.

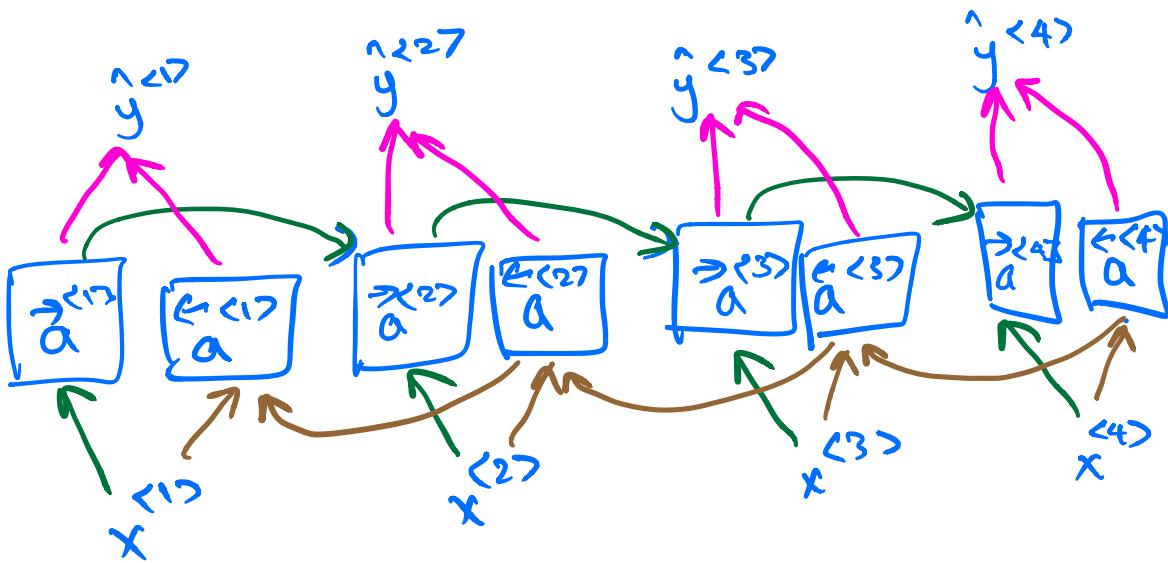
Variation:

Instead of just having gate values depend on $a^{(t-1)}, x^{(t)}$ we can add $c^{(t-1)}$. This is called peephole connection

GRU \rightarrow computationally effective

LSTM \rightarrow much effective

Bidirectional RNN:



acyclic
graph is
formed

→ step 1
→ step 2
→ step 3

$$\hat{y}^{(t)} = g(w_y [\vec{a}^{(t)}, \vec{a}^{(t)}] + b_y)$$

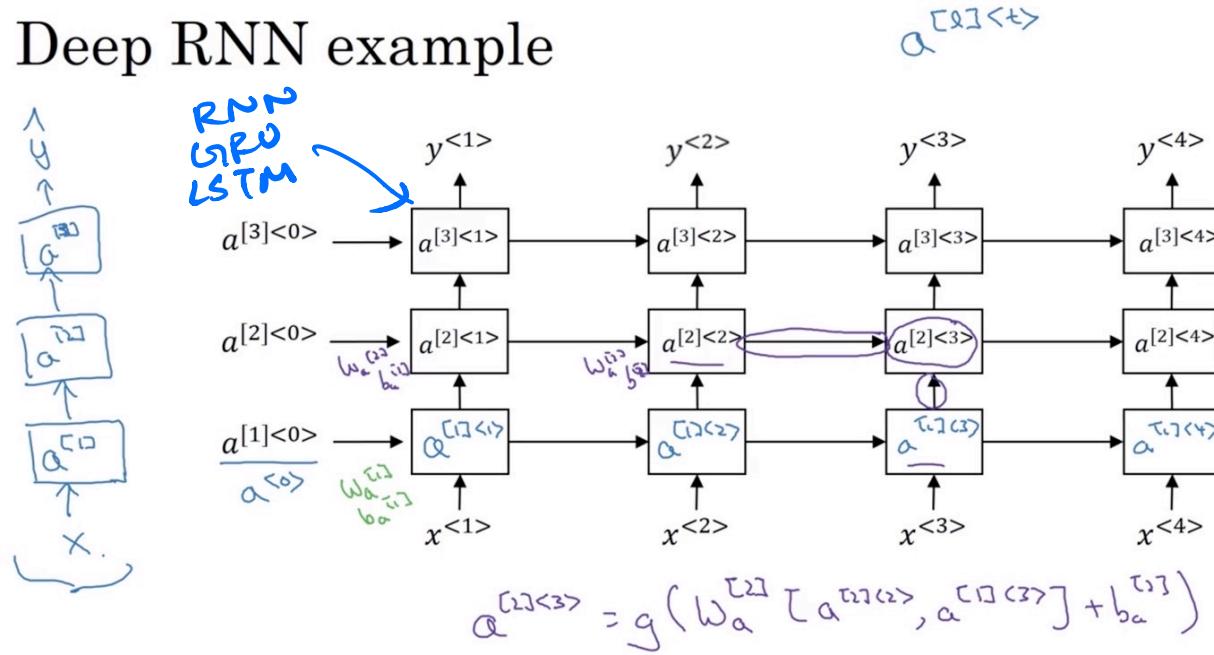
each box can be basic RNN
block or GRU or LSTM block

but we need entire sequence
before processing

Deep RNN:

even 3 layers stacked can take significant time.

Deep RNN example



we can also take $y^{<t>}$ used as input for further dnn without horizontal connection.

word representation:

$O_x \rightarrow$ one hot vector with
 x^{th} position as 1.

This doesn't know the relation
between words. Always O_x and

O_y will be zero.

featureized representation:
word embedding

This will learn values for
each feature for the word.
similar words will have
somewhat common features

Featurized representation: word embedding

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
300 Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
size cost alive verb	⋮ e_{5391}	⋮ e_{9853}				

I want a glass of orange juice.
I want a glass of apple juice.
Andrew Ng

using t-SNE to plot the embedding in 2D, we can see that more related words are closer together.

(e) orange and apple are closer than man and apple.

∴ we can generalize better in test data.

Transfer learning and word embeddings

1. Learn word embeddings from large text corpus. (1-100B words)

(Or download pre-trained embedding online.)

2. Transfer embedding to new task with smaller training set.
(say, 100k words) $\rightarrow 10,000 \rightarrow 300$

3. Optional: Continue to finetune the word embeddings with new data.

In face encoding \rightarrow we expect
to input unknown face to
network and it still gives
us encoding.

In word embedding \rightarrow we have
fixed vocabulary and
learn encoding for the
words.

Analogy

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{?}}$
 $e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}}$

$\frac{e_{\text{man}} - e_{\text{woman}}}{\|e_{\text{man}} - e_{\text{woman}}\|} \approx \frac{e_{\text{king}} - e_{\text{?}}}{\|e_{\text{king}} - e_{\text{?}}\|}$
 $\frac{e_{\text{man}} - e_{\text{woman}}}{\|e_{\text{man}} - e_{\text{woman}}\|} \approx \frac{e_{\text{king}} - e_{\text{queen}}}{\|e_{\text{king}} - e_{\text{queen}}\|}$

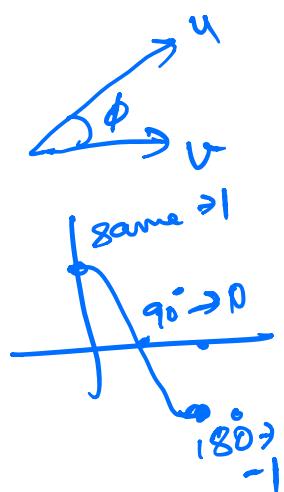
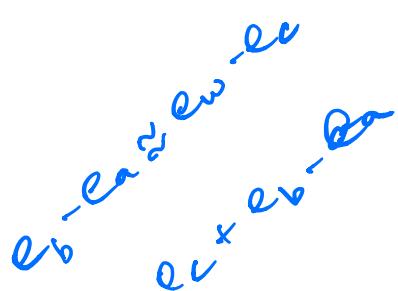
find word w :

$\underset{w}{\operatorname{argmax}}$ similarity (e_w ,
 $e_{\text{king}} - e_{\text{man}}$
 $+ e_{\text{woman}})$)

most common similarity:

cosine similarity:

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$



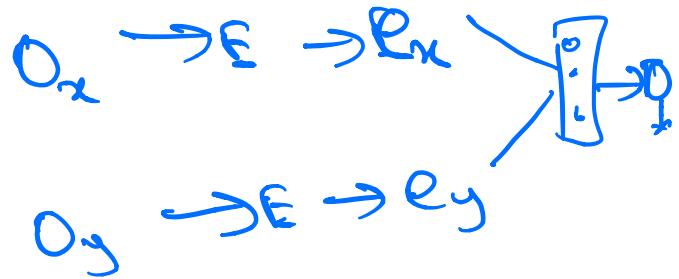
Neural language model: (to learn) (embedding)

i) construct one hot vector for each word

ii) $E \cdot O_x \rightarrow P_x$ for all words

iii) give O_x for necessary words as an input to a neural network that has softmax final layer to predict next word.

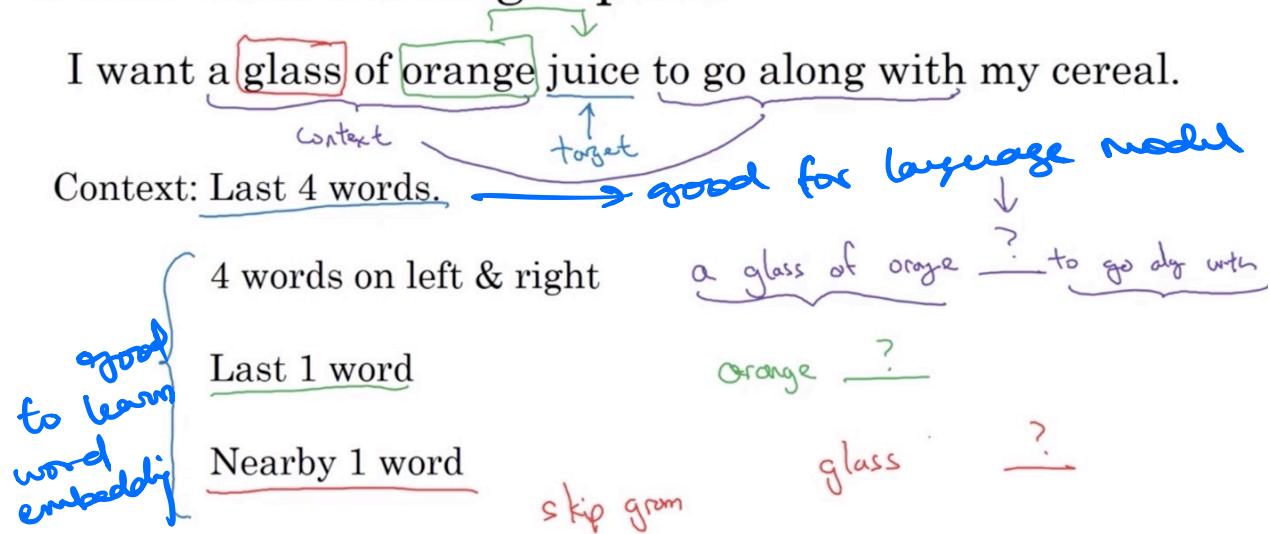
instead of all words we can also have fixed historical data like last 4 words



parameters $\rightarrow E$, weights of neural network

use gradient descent to find E .

Other context/target pairs



Skip gram:

* context - target pair for training → randomly chosen in a window

$x \rightarrow y$
 context $c \rightarrow$ target t

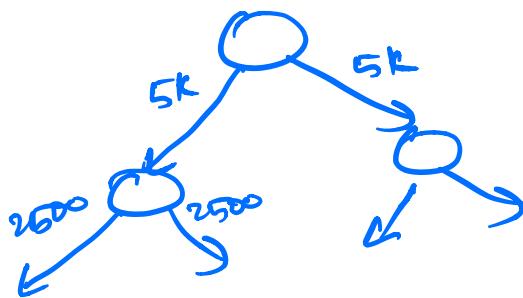
$$o_c \rightarrow E \rightarrow e_c \xrightarrow{\text{softmax}} \hat{y} \rightarrow \hat{y}$$

$$P(t|c) = \frac{e^{t^T e_c}}{\sum_j e^{j^T e_c}}$$

$$h(\hat{y}, y) = -\sum_{i=1}^{10000} y_i \log \hat{y}_i$$

The summation in $P(t|c)$ is computationally expensive

hierarchical softmax



https://www.youtube.com/watch?v=pzylWCelt_E

In practice, we have unbalanced tree to have common words are in higher level.

selecting $c \rightarrow$ uniformly random with some heuristics to select uncommon words.

selecting $t \rightarrow$ random word in a window.

New learning problem:

context - word
 x

is this
context
target pair

context \rightarrow randomly from sentence

target word \rightarrow { word in some, the
window, except

every the sample,
we have K - ve.

random word - ve
from dictionary, except

model:

x \rightarrow is target
(context, word) y
 \downarrow w

$$P(y=1 | c, t) = \sigma(\theta^T e_c)$$

$\theta^T e_c \rightarrow F \rightarrow e_x \rightarrow$ {
binary classifi

each word \rightarrow is target or not

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum f(w_i)^{3/4}}$$

Glove: global vector for word representation.

$X_{ij} = \# \text{ times } j \text{ appears in context of } i$

minimize :

$$\sum_{i=1}^{10000} \sum_{j=1}^{10000} f(x_{ij}) \left(\theta_i^T e_j + b_i + b_j^T e_c - \log x_{ij} \right)$$

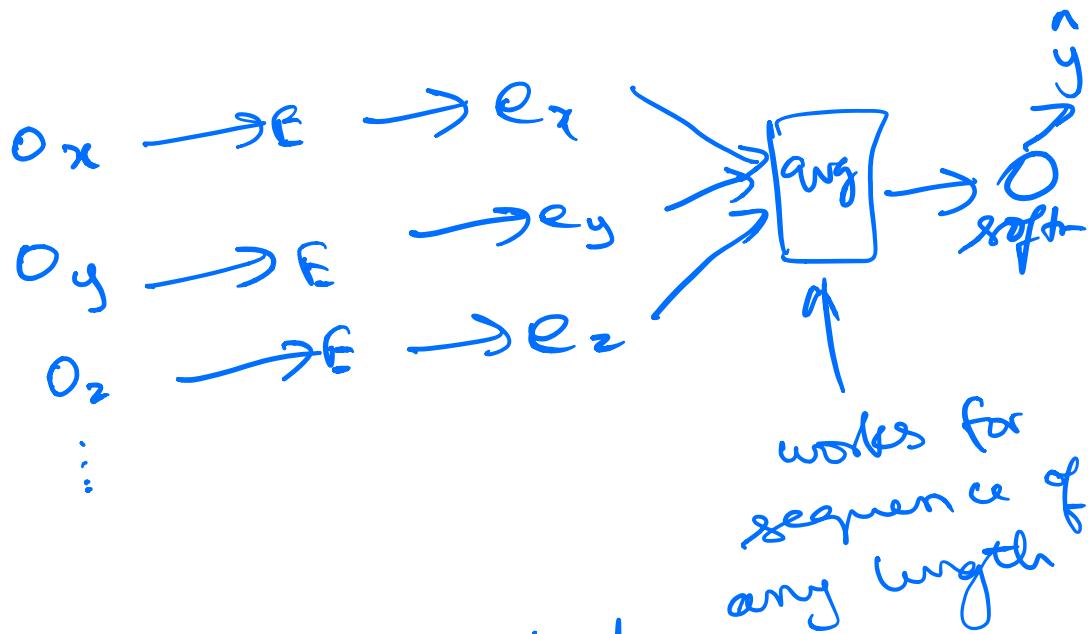
↑
weighting term

≥ 0 when $x_{ij} > 0$

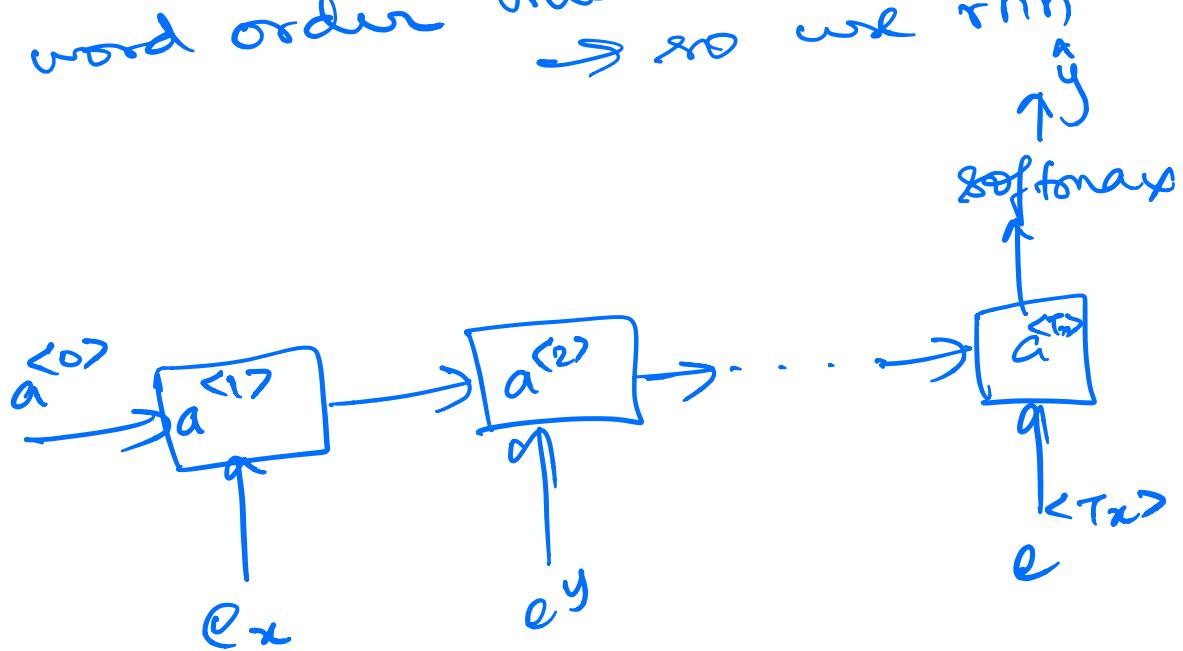
so that we won't have to calculate $\log 0$

Since θ_T and e_j are symmetric. So initialize randomly and $e_w^{(\text{final})} = \frac{e_w + \theta_w}{2}$

Sentiment classification:

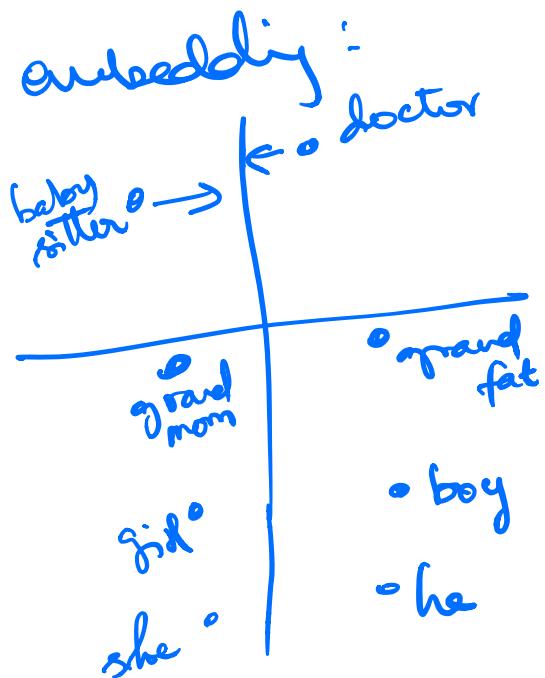


But this doesn't take word order into consideration
→ so we use RNN



Debiasing word embedding:

man : programmer :: women :
home maker
(⊗ wrong)



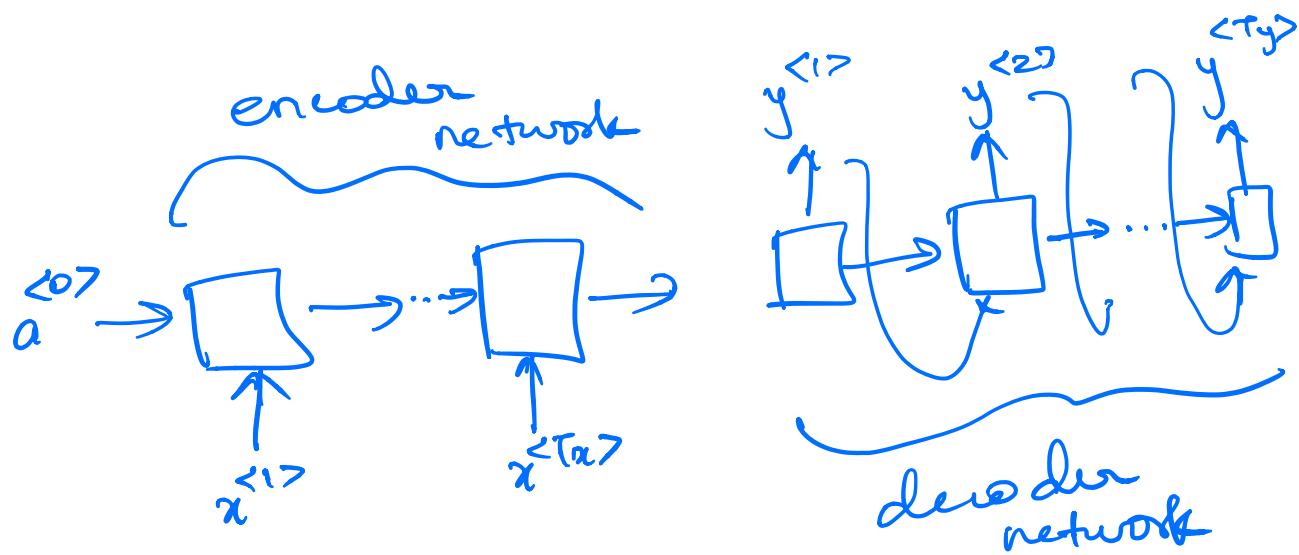
to maintain
distance between
the "boy - girl"
and "grand mom - fat"
same

1. identify bias direction
She - like

2. neutralize:
for word that
is not definit
(doctor),
project to get
rid of bias

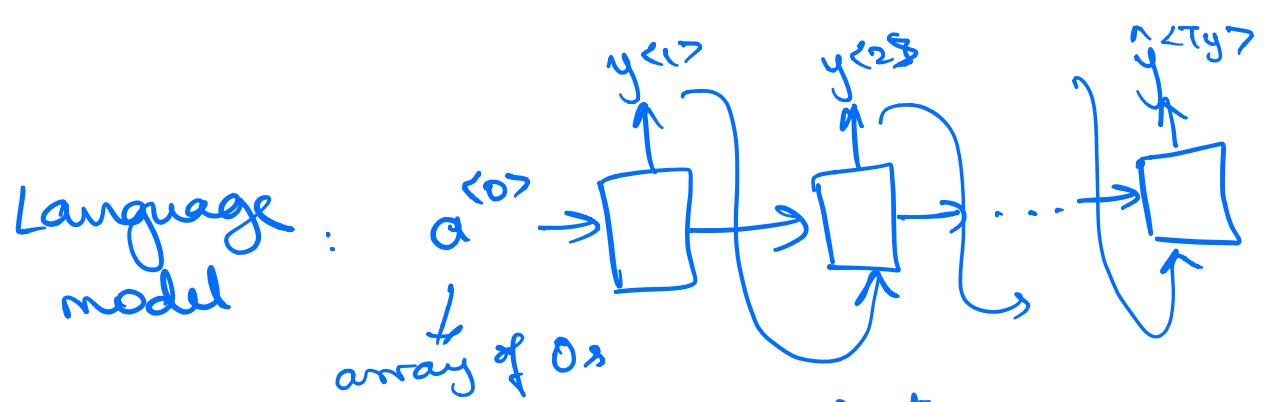
3. equalize pairs
(boy - girl)
so move boy +
girl to be
same distance
to that word

Sequence to sequence model:



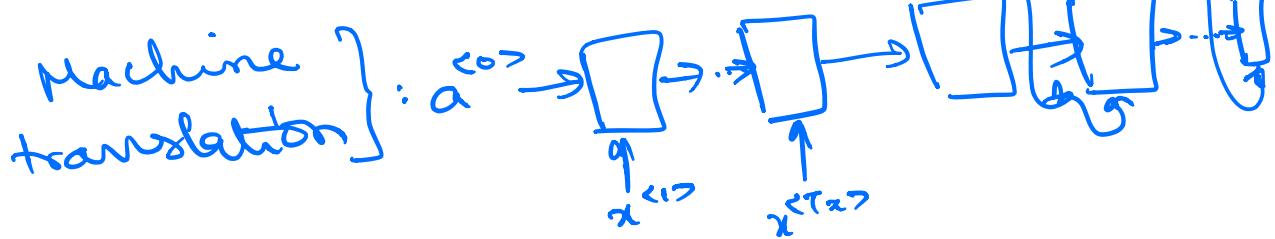
works well when
enough dataset is given.

For image captioning, we a
pretrained CNN as encoder
network



this calculates

$$P(y^{<1>} \dots y^{<Ty>})$$



here the decoder part is
same as the one we used

for language model.

Except instead of $a^{<0>} = \vec{0}$,
here we have $\vec{q}^{<0>}$ as encoder
network's output.

hence also
called
conditional
language
model

∴ here we calculate

$$P(\hat{y}^{<1>} \dots \hat{y}^{<Ty>} | x^{<1>} \dots x^{<Tx>})$$

for french \rightarrow english

$$P(y^{<1>} \dots y^{<Ty>} | x)$$

↑ ↑
english french

we will have probabilities for diff.
english sentences.

(we cannot sample randomly
from that)

so find $y^{<1>} \dots y^{<Ty>}$ that maximizes

$$P(y|x)$$

\Rightarrow we use beam search

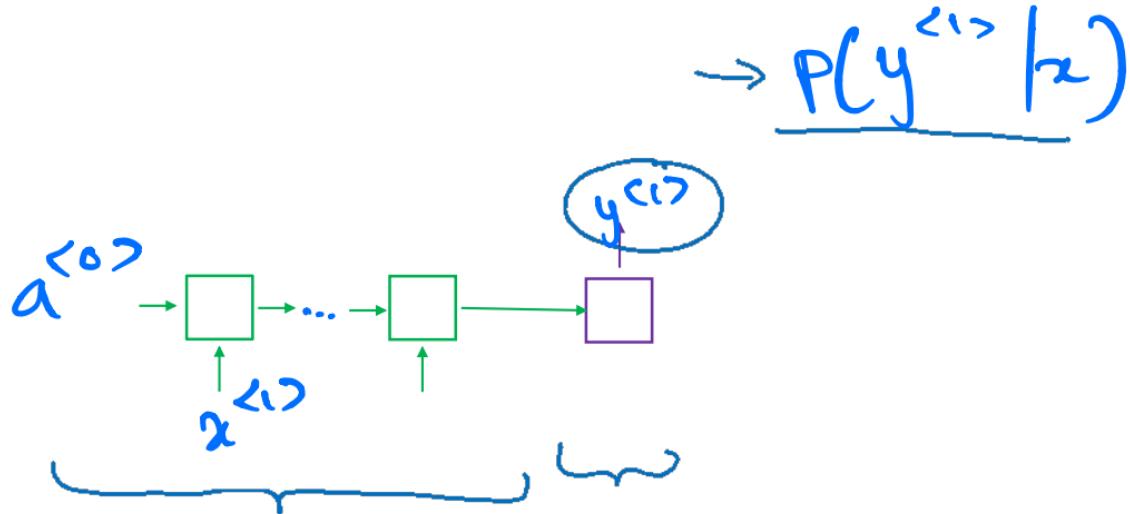
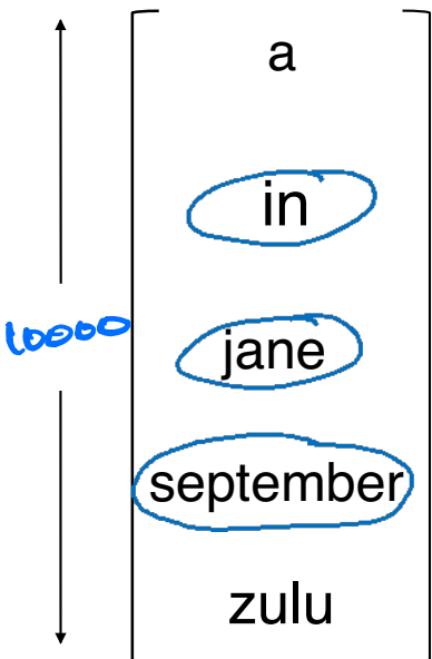
greedy search (each word selected
will have max prob)

\Rightarrow overall optimality may not
be maintained.

Beam search algorithm

B = 3 (beam width)

Step 1



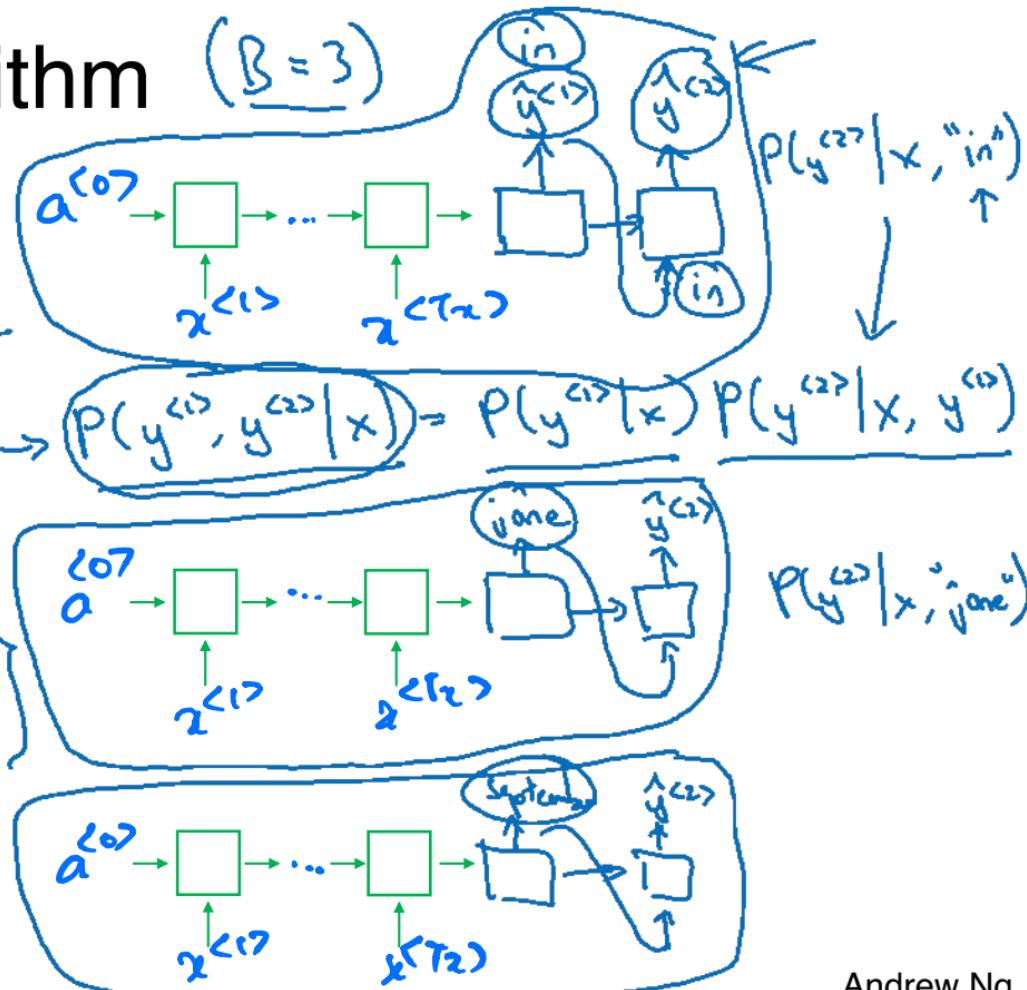
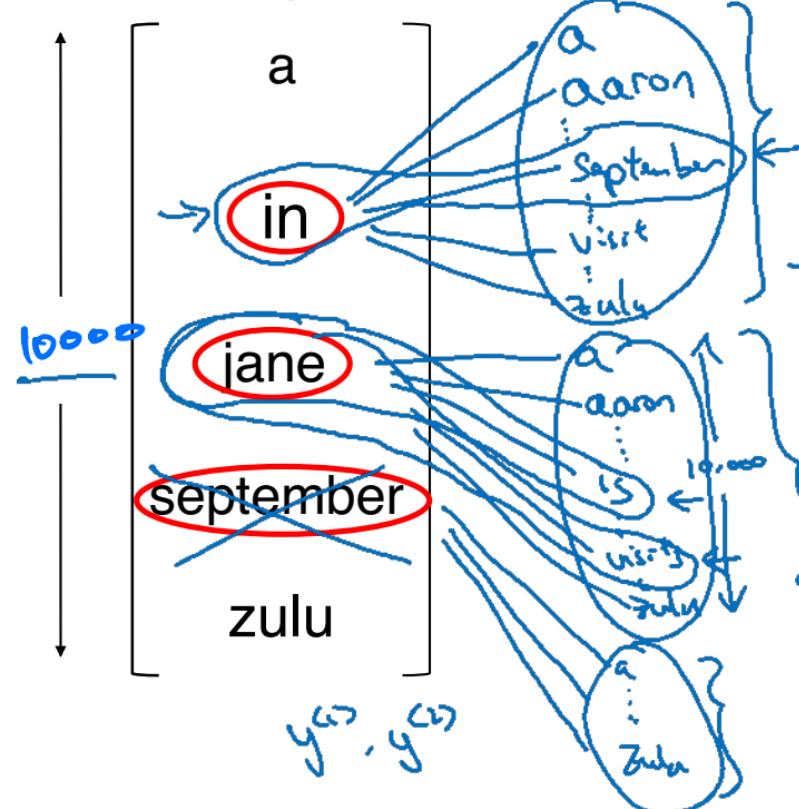
B copies of network
needed

Beam search algorithm

$(B = 3)$

Step 1

Step 2



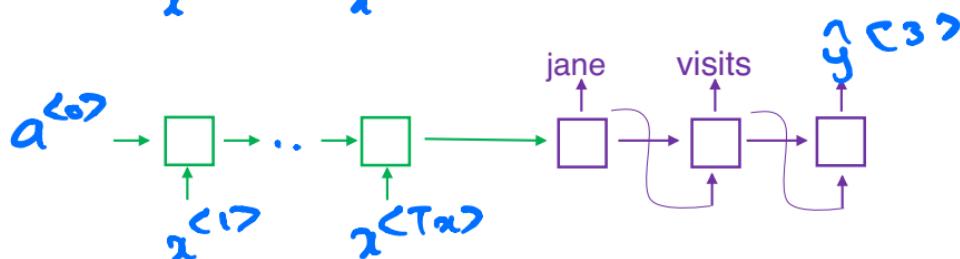
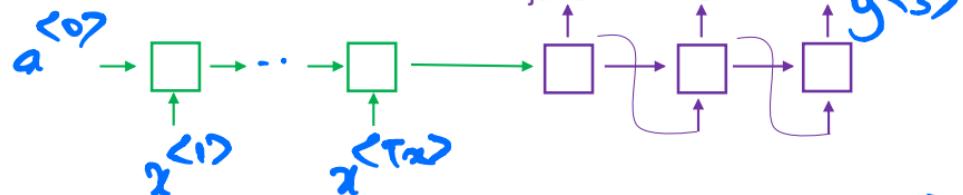
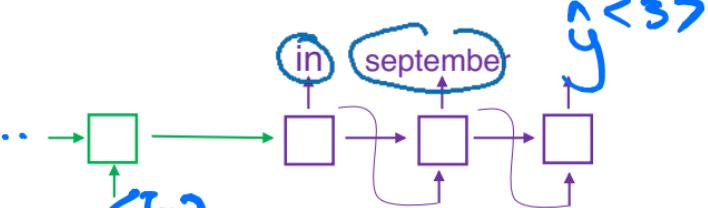
Beam search ($B=3$)

in september
aaron
jane
zulu

jane is
visits
zulu

jane visits
africa
zulu

$B=1 \rightarrow$ greedy search



jane visits africa in september. <EOS>

length normalization:

we have :

$$\arg \max_y \prod_{t=1}^{Ty} P(y^{<t>} | x, y^{<1>}..y^{<t-1>})$$

as these are small values,
we may have underflow.

$$\therefore \arg \max_y \sum_{t=1}^{Ty} \log(P(y^{<t>} | x, y^{<1>}..y^{<t-1>}))$$

This tends to prefer short output

so normalize it:

$$\arg \max_y \frac{1}{Ty^\alpha} \sum_{t=1}^{Ty} \log(P(y^{<t>} | x, y^{<1>}..y^{<t-1>}))$$

$\alpha = 1 \rightarrow$ full normalization

$\alpha = 0 \rightarrow$ no normalization

large B : better result (check many possib)
slower

small B : worse result
faster

unlike BFS/DFS, beam search runs faster but is not guaranteed to find exact maximum.

Error analysis:

statement, y^* , \hat{y}
 \uparrow \uparrow
human RNN + beam search

case 1: $P(y^* | x) > P(\hat{y} | x)$

\Rightarrow beam search fault

case 2: $P(y^* | x) \leq P(\hat{y} | x)$

\Rightarrow RNN fault

BLEU score: bilingual evaluation understudy

how good is the machine translation

If predicted statement ≈ reference statement
→ high bleu score

Give credit upto max. no. of times words occur in reference sentence

Bleu score on bigrams

Example: Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: The cat the cat on the mat. ←

	Count	Count clip	
the cat	2 ←	1 ←	
cat the	1 ←	0	
cat on	1 ←	1 ←	
on the	1 ←	1 ←	
the mat	1 ←	1 ←	

$\frac{4}{6}$.

Bleu score on unigrams

Example: Reference 1: The cat is on the mat.

$$P_1, P_2, = 1.0$$

Reference 2: There is a cat on the mat.

→ MT output: The cat the cat on the mat. (\hat{y})

$$P_1 = \frac{\sum_{\text{Unigrams } \in \hat{y}} \text{Count}_{clip}(\text{unigram})}{\sum_{\text{Unigrams } \in \hat{y}} \text{Count}(\text{unigram})}$$

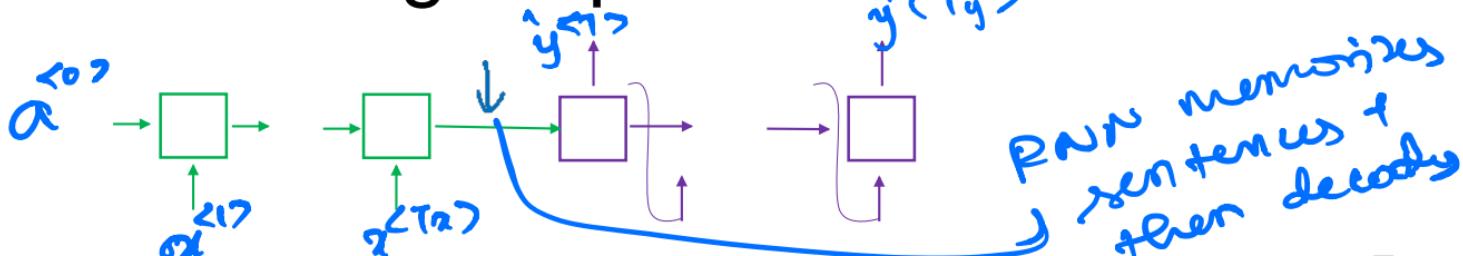
$$P_n = \frac{\sum_{n\text{-grams } \in \hat{y}} \text{Count}_{clip}(n\text{-gram})}{\sum_{n\text{-grams } \in \hat{y}} \text{Count}(n\text{-gram})}$$

combined blue score } = BP \exp \left(\frac{1}{4} \sum_{n=1}^4 P_n \right)

brevity penalty → to penalize short sentence

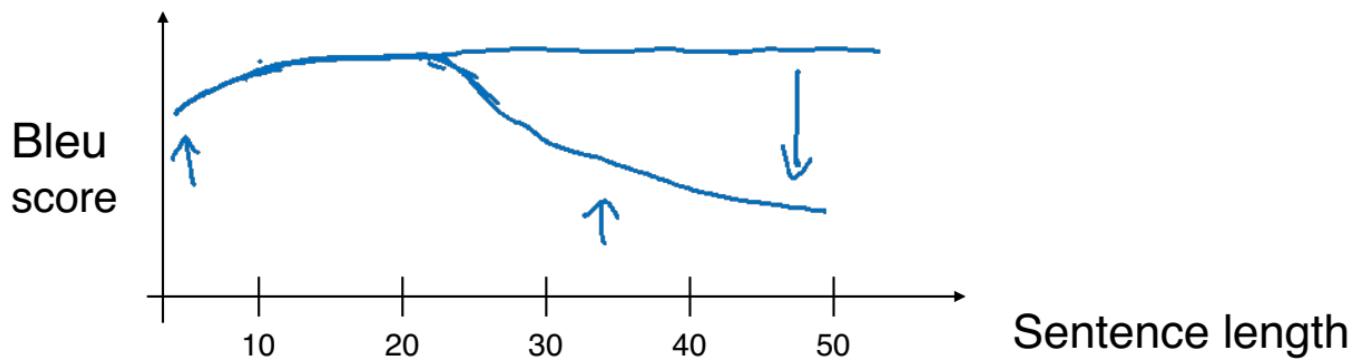
$$BP = \begin{cases} 1, & \text{if } \text{MT-out-len} > \text{reference-out-l} \\ \exp(1 - \frac{\text{ref-o-len}}{\text{MT-o-len}}), & \text{otherwise} \end{cases}$$

The problem of long sequences

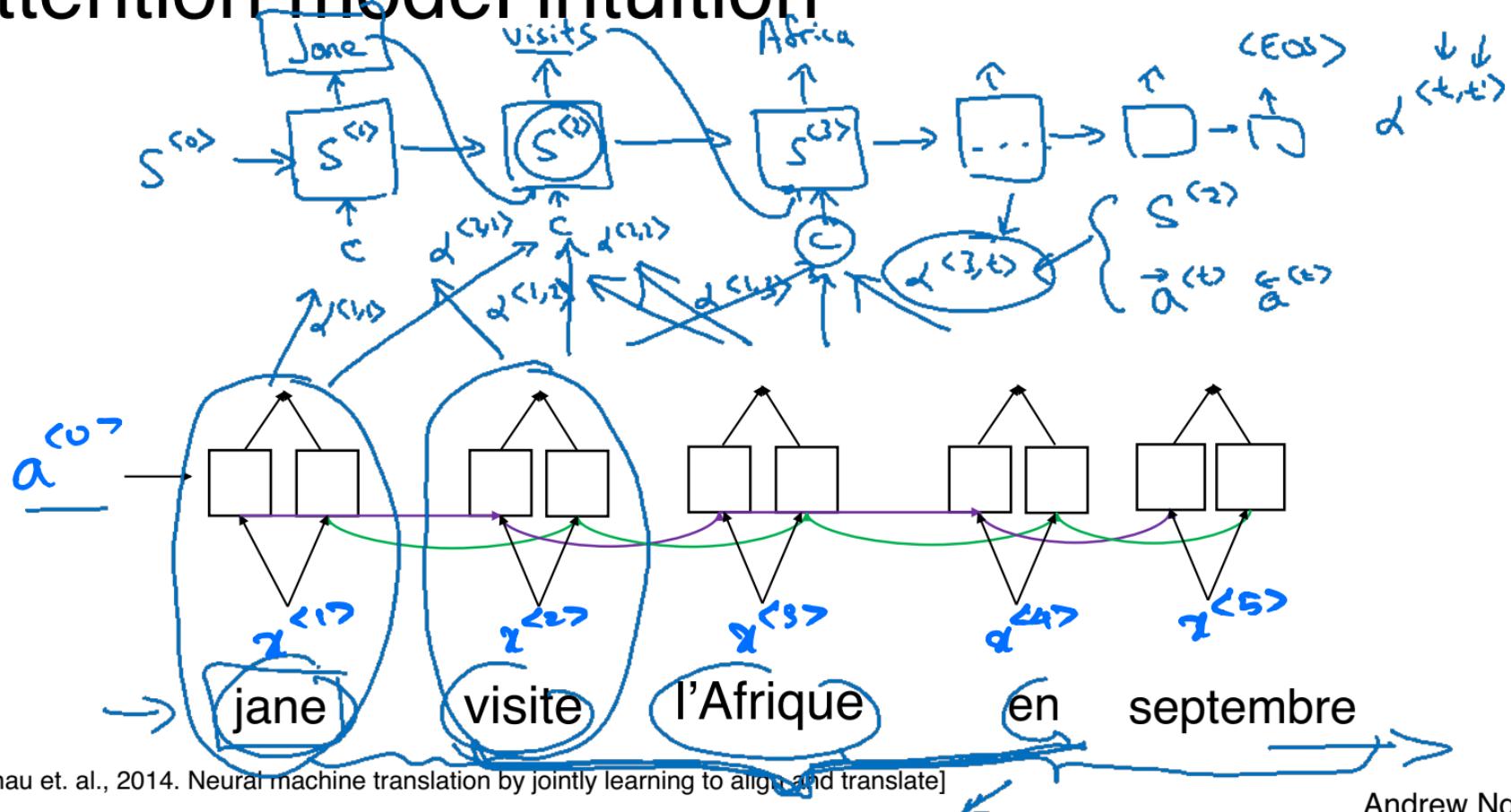


Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

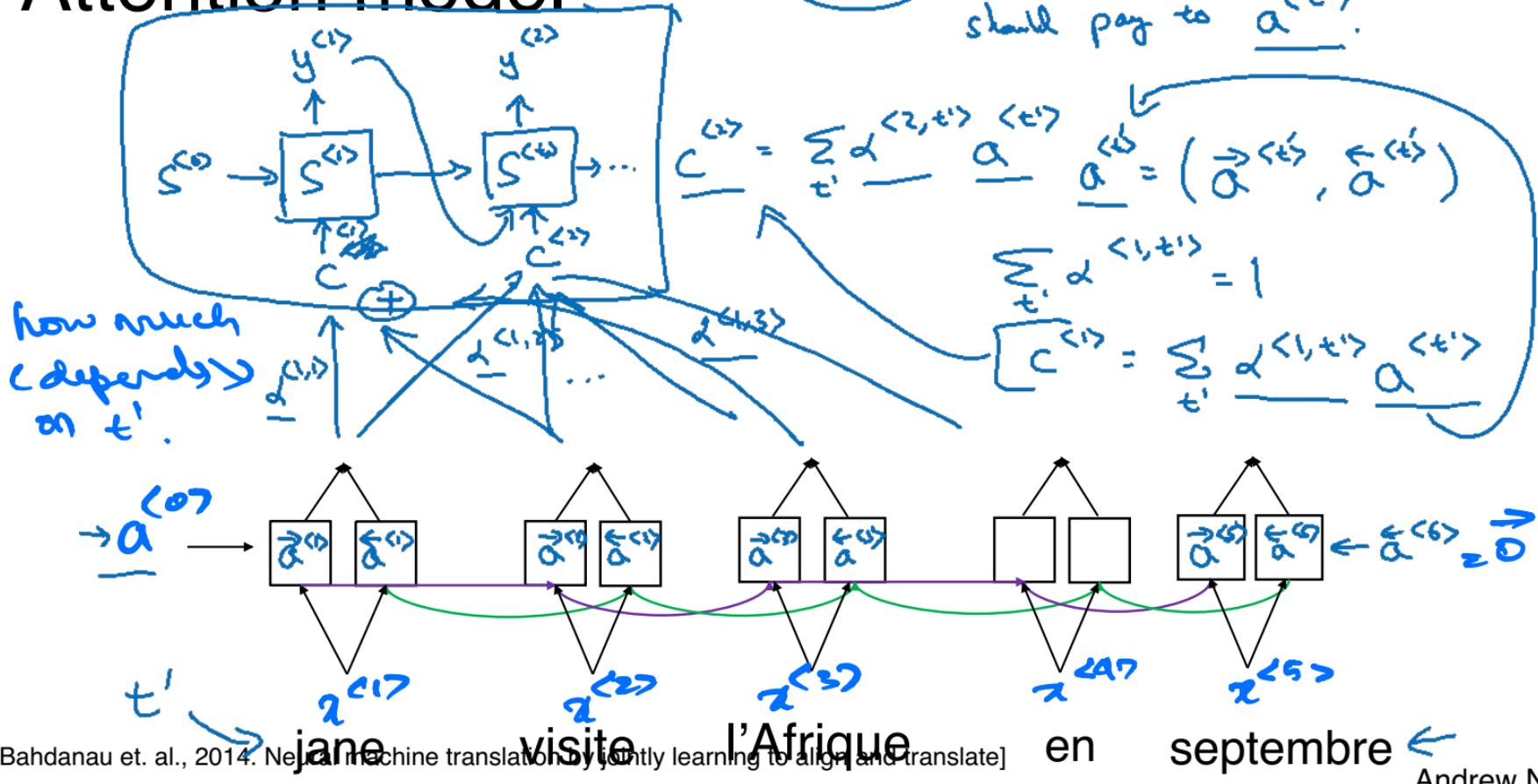
Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go too.



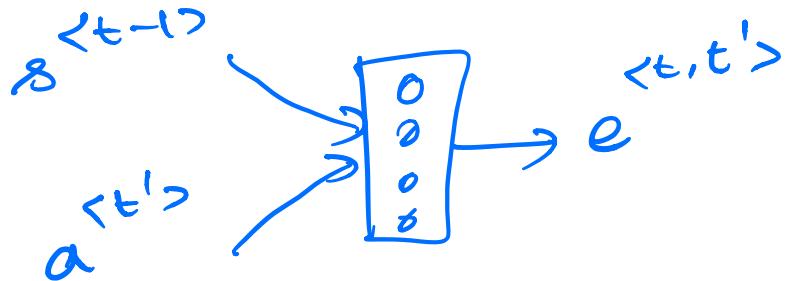
Attention model intuition



Attention model



$$\alpha^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t'=1}^{T_x} \exp(e^{(t,t')})}$$



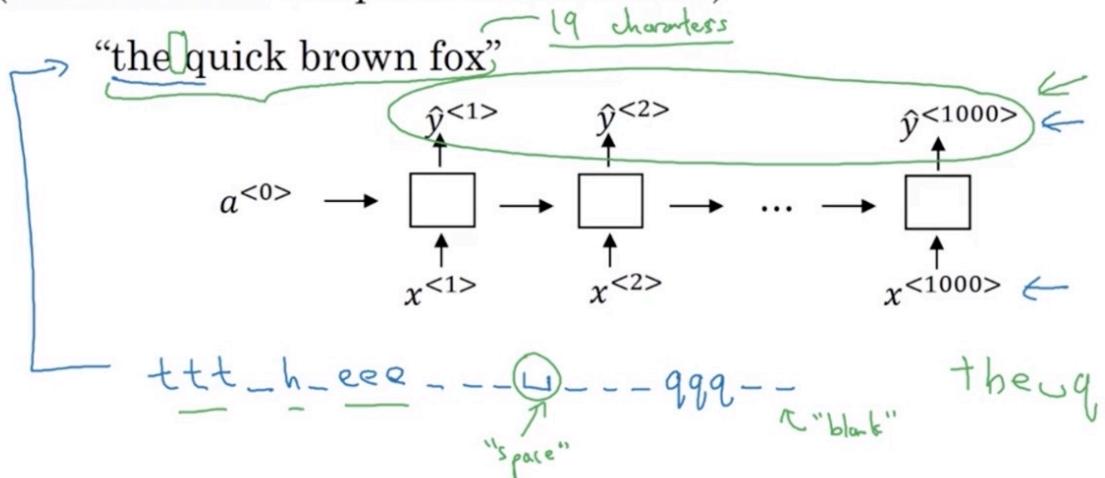
Time complexity: $T_x T_y$

spectrogram
→ common
audio pre
processing

Trigger word detection:
train RNN to output 0 when
user hasn't said trigger word
and output 1 when user
says the trigger word.

CTC cost for speech recognition

(Connectionist temporal classification)



Basic rule: collapse repeated characters not separated by “blank”